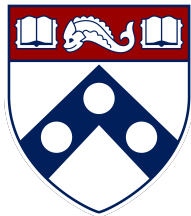


# Non-Asymptotic Analysis for Reinforcement Learning (Part 2)



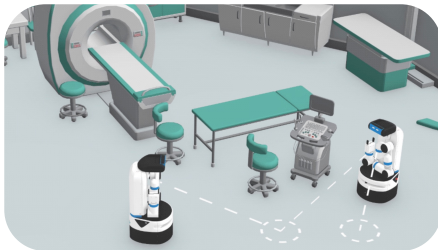
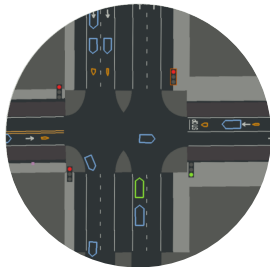
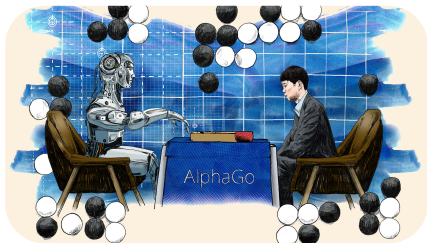
Yuxin Chen

Wharton Statistics & Data Science, SIGMETRICS 2023

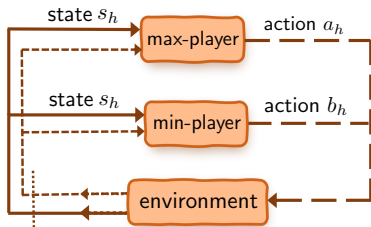
*Multi-agent RL with a generative model*

# Multi-agent reinforcement learning (MARL)

---

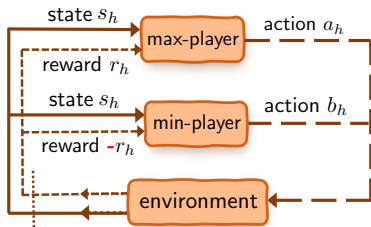


# Two-player zero-sum Markov games (finite-horizon)



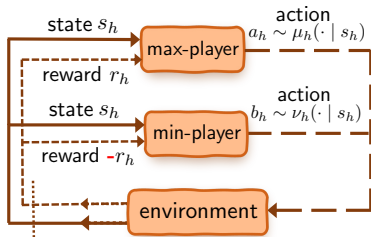
- $\mathcal{S} = [S]$ : state space
- $\mathcal{A} = [A]$ : action space of max-player
- $H$ : horizon
- $\mathcal{B} = [B]$ : action space of min-player

# Two-player zero-sum Markov games (finite-horizon)



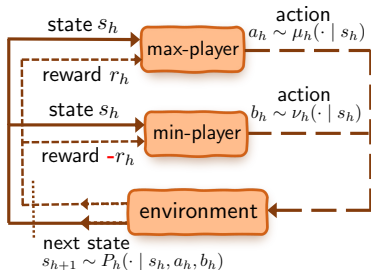
- $\mathcal{S} = [S]$ : state space
- $\mathcal{A} = [A]$ : action space of max-player
- $H$ : horizon
- $\mathcal{B} = [B]$ : action space of min-player
- immediate reward: max-player  $r(s, a, b) \in [0, 1]$   
min-player  $-r(s, a, b)$

# Two-player zero-sum Markov games (finite-horizon)



- $\mathcal{S} = [S]$ : state space
- $\mathcal{A} = [A]$ : action space of max-player
- $H$ : horizon
- $\mathcal{B} = [B]$ : action space of min-player
- immediate reward: max-player  $r(s, a, b) \in [0, 1]$   
min-player  $-r(s, a, b)$
- $\mu : \mathcal{S} \times [H] \rightarrow \Delta(\mathcal{A})$ : policy of max-player
- $\nu : \mathcal{S} \times [H] \rightarrow \Delta(\mathcal{B})$ : policy of min-player

# Two-player zero-sum Markov games (finite-horizon)



- $\mathcal{S} = [S]$ : state space
- $\mathcal{A} = [A]$ : action space of max-player
- $H$ : horizon
- $\mathcal{B} = [B]$ : action space of min-player
- immediate reward: max-player  $r(s, a, b) \in [0, 1]$   
min-player  $-r(s, a, b)$
- $\mu : \mathcal{S} \times [H] \rightarrow \Delta(\mathcal{A})$ : policy of max-player  
 $\nu : \mathcal{S} \times [H] \rightarrow \Delta(\mathcal{B})$ : policy of min-player
- $P_h(\cdot | s, a, b)$ : **unknown** transition probabilities

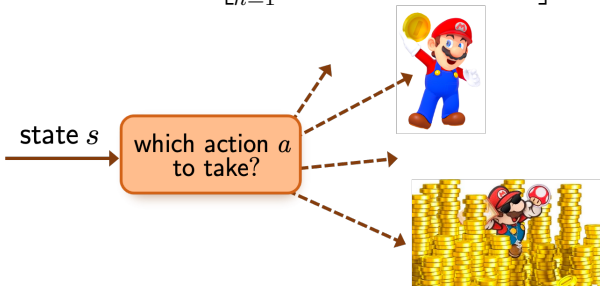
**Value function** under *independent* policies  $(\mu, \nu)$  (no coordination)

$$V^{\mu, \nu}(s) := \mathbb{E} \left[ \sum_{h=1}^H r_h(s_h, a_h, b_h) \mid s_1 = s \right]$$



**Value function** under *independent* policies  $(\mu, \nu)$  (no coordination)

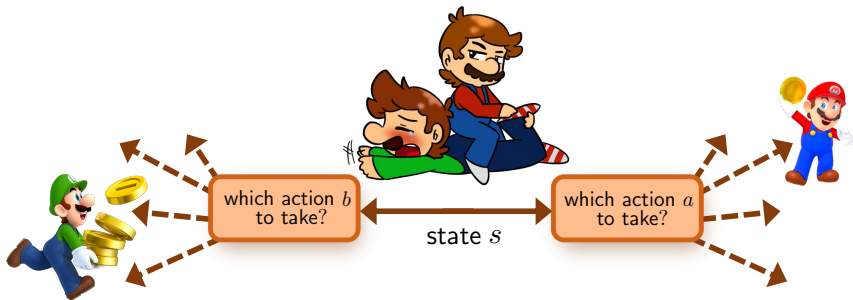
$$V^{\mu, \nu}(s) := \mathbb{E} \left[ \sum_{h=1}^H r_h(s_h, a_h, b_h) \mid s_1 = s \right]$$



- Each agent seeks **optimal policy** maximizing her own value

## Value function under *independent* policies $(\mu, \nu)$ (no coordination)

$$V^{\mu, \nu}(s) := \mathbb{E} \left[ \sum_{h=1}^H r_h(s_h, a_h, b_h) \mid s_1 = s \right]$$



- Each agent seeks **optimal policy** maximizing her own value
- But two agents have conflicting goals ...

# Compromise: Nash equilibrium (NE)

---



*John von Neumann*



*John Nash*

An NE policy pair  $(\mu^*, \nu^*)$  obeys

$$\max_{\mu} V^{\mu, \nu^*} = V^{\mu^*, \nu^*} = \min_{\nu} V^{\mu^*, \nu}$$

# Compromise: Nash equilibrium (NE)

---



*John von Neumann*



*John Nash*

An NE policy pair  $(\mu^*, \nu^*)$  obeys

$$\max_{\mu} V^{\mu, \nu^*} = V^{\mu^*, \nu^*} = \min_{\nu} V^{\mu^*, \nu}$$

- no unilateral deviation is beneficial

# Compromise: Nash equilibrium (NE)

---



*John von Neumann*



*John Nash*

An NE policy pair  $(\mu^*, \nu^*)$  obeys

$$\max_{\mu} V^{\mu, \nu^*} = V^{\mu^*, \nu^*} = \min_{\nu} V^{\mu^*, \nu}$$

- no unilateral deviation is beneficial
- no coordination between two agents (they act *independently*)

# Compromise: Nash equilibrium (NE)

---



*John von Neumann*



*John Nash*

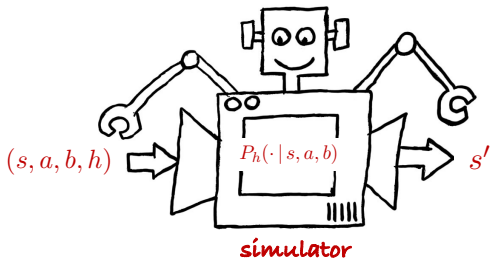
An  $\varepsilon$ -NE policy pair  $(\hat{\mu}, \hat{\nu})$  obeys

$$\max_{\mu} V^{\mu, \hat{\nu}} - \varepsilon \leq V^{\hat{\mu}, \hat{\nu}} \leq \min_{\nu} V^{\hat{\mu}, \nu} + \varepsilon$$

- no unilateral deviation is beneficial
- no coordination between two agents (they act *independently*)

# Learning NEs with a simulator

---

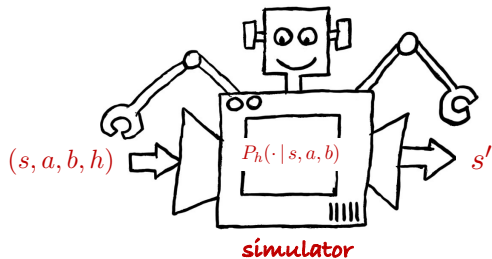


**input:** any  $(s, a, b, h)$

**output:** an independent sample  $s' \sim P_h(\cdot | s, a, b)$

# Learning NEs with a simulator

---



**input:** any  $(s, a, b, h)$

**output:** an independent sample  $s' \sim P_h(\cdot | s, a, b)$

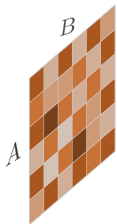
**Question:** how many samples are sufficient to learn an  $\varepsilon$ -Nash policy pair?



# Model-based approach (non-adaptive sampling)

---

— Zhang, Kakade, Başar, Yang '20

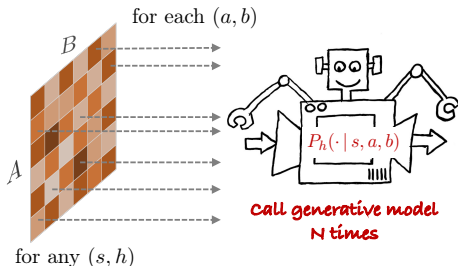


for any  $(s, h)$

1. for each  $(s, a, b, h)$ , call simulator  $N$  times

# Model-based approach (non-adaptive sampling)

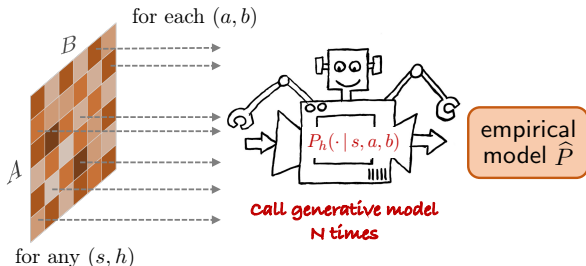
— Zhang, Kakade, Başar, Yang '20



1. for each  $(s, a, b, h)$ , call simulator  $N$  times

# Model-based approach (non-adaptive sampling)

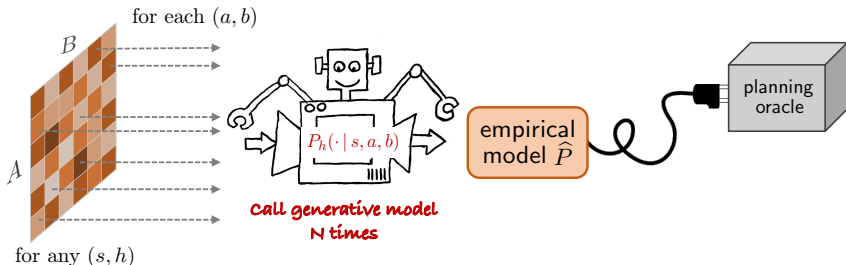
— Zhang, Kakade, Başar, Yang '20



1. for each  $(s, a, b, h)$ , call simulator  $N$  times
2. build empirical model  $\hat{P}$

# Model-based approach (non-adaptive sampling)

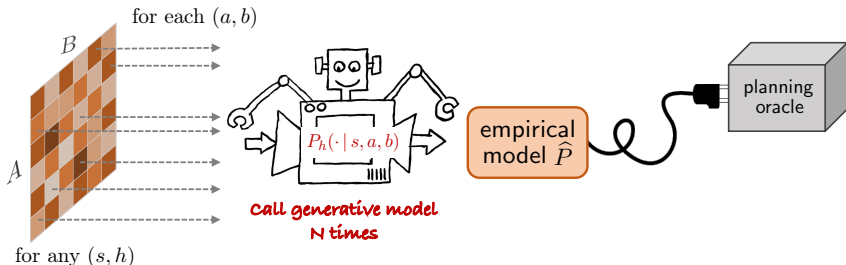
— Zhang, Kakade, Başar, Yang '20



1. for each  $(s, a, b, h)$ , call simulator  $N$  times
2. build empirical model  $\hat{P}$ , and run "plug-in" methods

# Model-based approach (non-adaptive sampling)

— Zhang, Kakade, Başar, Yang '20



1. for each  $(s, a, b, h)$ , call simulator  $N$  times
2. build empirical model  $\hat{P}$ , and run “plug-in” methods

sample complexity:  $\frac{H^4 SAB}{\epsilon^2}$

# Curse of multiple agents

---

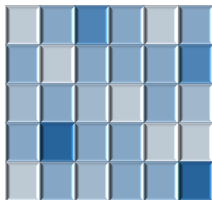


1 player:  $A$

Let's look at the **size** of joint action space ...

# Curse of multiple agents

---



1 player:  $A$



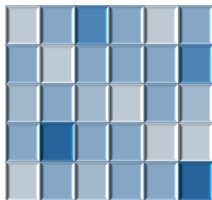
2 players:  $AB$

Let's look at the **size** of joint action space ...

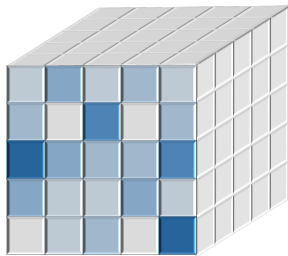
# Curse of multiple agents



1 player:  $A$



2 players:  $AB$



$m$  players:  $A_1A_2 \cdots A_m$

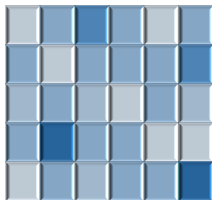
Let's look at the **size** of joint action space ...



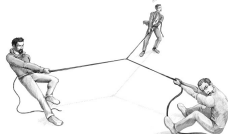
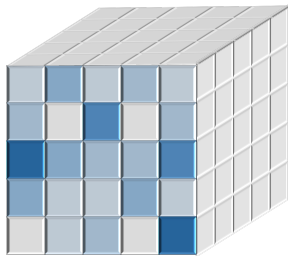
# Curse of multiple agents



1 player:  $A$

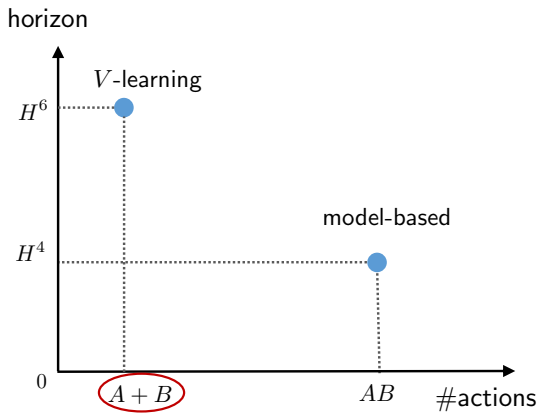


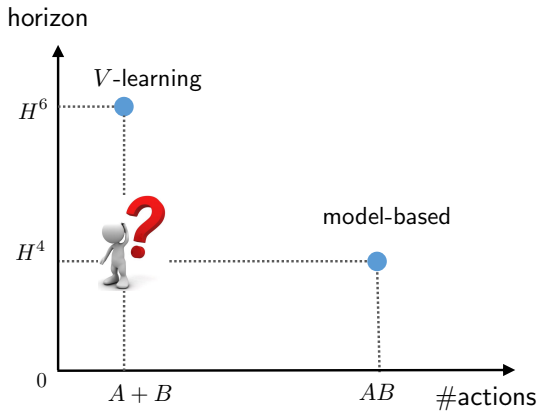
2 players:  $AB$

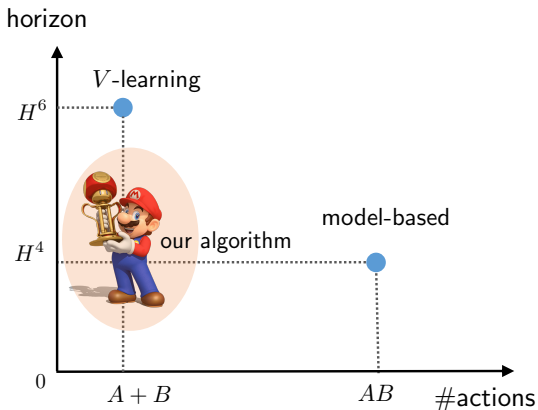


$m$  players:  $A_1A_2 \cdots A_m$

# joint actions **blows up** geometrically in # players!







### Theorem 1 (Li, Chi, Wei, Chen '22)

For any  $0 < \epsilon \leq H$ , one can design an algorithm that finds an  $\epsilon$ -Nash policy pair  $(\hat{\mu}, \hat{\nu})$  with high prob., with sample complexity at most

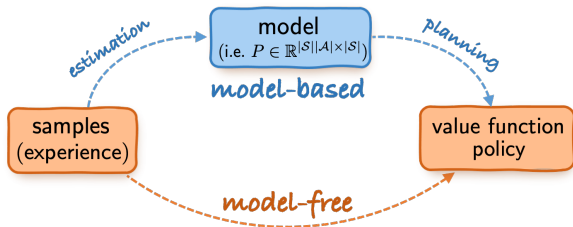
$$\tilde{O}\left(\frac{H^4 S(A+B)}{\epsilon^2}\right) \quad (\text{minimax-optimal } \forall \epsilon)$$

## Model-free / value-based RL

1. Basics of Q-learning
2. Synchronous Q-learning and variance reduction (simulator)
3. Asynchronous Q-learning (Markovian data)
4. Q-learning with lower confidence bounds (offline RL)
5. Q-learning with upper confidence bounds (online RL)

# Model-based vs. model-free RL

---

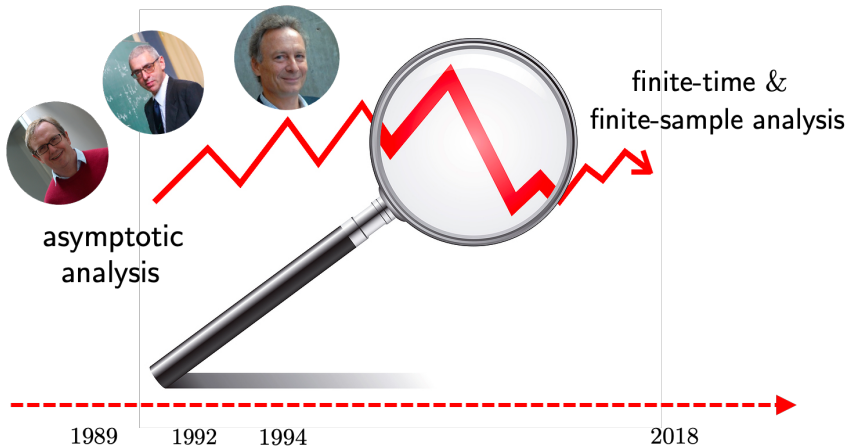


## Model-based approach (“plug-in”)

1. build empirical estimate  $\hat{P}$  for  $P$
2. planning based on empirical  $\hat{P}$

## Model-free / value-based approach

- learning w/o modeling & estimating environment explicitly
- memory-efficient, online, ...



Focus of this part: classical **Q-learning** algorithm and its variants

# A starting point: Bellman optimality principle

---

## Bellman operator

$$\mathcal{T}(Q)(s, a) := \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \underbrace{\max_{a' \in \mathcal{A}} Q(s', a')}_{\text{next state's value}} \right]$$

- one-step look-ahead



# A starting point: Bellman optimality principle

---

## Bellman operator

$$\mathcal{T}(Q)(s, a) := \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \underbrace{\max_{a' \in \mathcal{A}} Q(s', a')}_{\text{next state's value}} \right]$$

- one-step look-ahead

**Bellman equation:**  $Q^*$  is *unique* solution to

$$\mathcal{T}(Q^*) = Q^*$$

# A starting point: Bellman optimality principle

## Bellman operator

$$\mathcal{T}(Q)(s, a) := \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \underbrace{\max_{a' \in \mathcal{A}} Q(s', a')}_{\text{next state's value}} \right]$$

- one-step look-ahead

**Bellman equation:**  $Q^*$  is *unique* solution to

$$\mathcal{T}(Q^*) = Q^*$$

- **takeaway message:** it suffices to solve the Bellman equation
- **challenge:** how to solve it using stochastic samples?



*Richard Bellman*

# Q-learning: a stochastic approximation algorithm

---



Chris Watkins



Peter Dayan

Stochastic approximation for solving the **Bellman equation**

Robbins & Monro, 1951

$$\mathcal{T}(Q) - Q = 0$$

where

$$\mathcal{T}(Q)(s, a) := \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \underbrace{\max_{a' \in \mathcal{A}} Q(s', a')}_{\text{next state's value}} \right].$$

# Q-learning: a stochastic approximation algorithm

---



*Chris Watkins*



*Peter Dayan*

Stochastic approximation for solving Bellman equation  $\mathcal{T}(Q) - Q = 0$

$$\underbrace{Q_{t+1}(s, a) = Q_t(s, a) + \eta_t(\mathcal{T}_t(Q_t)(s, a) - Q_t(s, a))}_{\text{sample transition } (s, a, s')}, \quad t \geq 0$$

# Q-learning: a stochastic approximation algorithm

---



Chris Watkins



Peter Dayan

Stochastic approximation for solving Bellman equation  $\mathcal{T}(Q) - Q = 0$

$$\underbrace{Q_{t+1}(s, a) = Q_t(s, a) + \eta_t(\mathcal{T}_t(Q_t)(s, a) - Q_t(s, a))}_{\text{sample transition } (s, a, s')}, \quad t \geq 0$$

$$\mathcal{T}_t(Q)(s, a) = r(s, a) + \gamma \max_{a'} Q(s', a')$$

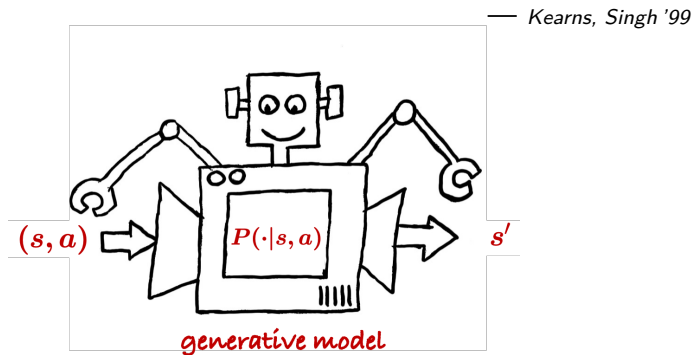
$$\mathcal{T}(Q)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \max_{a'} Q(s', a') \right]$$

## Model-free RL

1. Basics of Q-learning
2. Synchronous Q-learning and variance reduction (simulator)
3. Asynchronous Q-learning (Markovian data)
4. Q-learning with lower confidence bounds (offline RL)
5. Q-learning with upper confidence bounds (online RL)

# A generative model / simulator

---



Each iteration, draw an independent sample  $(s, a, s')$  for given  $(s, a)$

# Synchronous Q-learning

---



Chris Watkins



Peter Dayan

**for**  $t = 0, 1, \dots, T$

**for** each  $(s, a) \in \mathcal{S} \times \mathcal{A}$

draw a sample  $(s, a, s')$ , run

$$Q_{t+1}(s, a) = (1 - \eta_t)Q_t(s, a) + \eta_t \left\{ r(s, a) + \gamma \max_{a'} Q_t(s', a') \right\}$$

**synchronous:** all state-action pairs are updated simultaneously

- total sample size:  $T|\mathcal{S}||\mathcal{A}|$



# Sample complexity of synchronous Q-learning

## Theorem 2 (Li, Cai, Chen, Wei, Chi '21)

For any  $0 < \varepsilon \leq 1$ , synchronous Q-learning yields  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  with high prob. and  $\mathbb{E}[\|\hat{Q} - Q^*\|_\infty] \leq \varepsilon$ , with sample size **at most**

$$\begin{cases} \tilde{O}\left(\frac{|S||\mathcal{A}|}{(1-\gamma)^4 \varepsilon^2}\right) & \text{if } |\mathcal{A}| \geq 2 \\ \tilde{O}\left(\frac{|S|}{(1-\gamma)^3 \varepsilon^2}\right) & \text{if } |\mathcal{A}| = 1 \end{cases} \quad (\text{TD learning})$$

# Sample complexity of synchronous Q-learning

## Theorem 2 (Li, Cai, Chen, Wei, Chi '21)

For any  $0 < \varepsilon \leq 1$ , synchronous Q-learning yields  $\|\widehat{Q} - Q^*\|_\infty \leq \varepsilon$  with high prob. and  $\mathbb{E}[\|\widehat{Q} - Q^*\|_\infty] \leq \varepsilon$ , with sample size **at most**

$$\begin{cases} \tilde{O}\left(\frac{|S||\mathcal{A}|}{(1-\gamma)^4 \varepsilon^2}\right) & \text{if } |\mathcal{A}| \geq 2 \\ \tilde{O}\left(\frac{|S|}{(1-\gamma)^3 \varepsilon^2}\right) & \text{if } |\mathcal{A}| = 1 \end{cases} \quad (\text{TD learning})$$

- Covers both *constant* and *rescaled linear* learning rates:

$$\eta_t \equiv \frac{1}{1 + \frac{c_1(1-\gamma)T}{\log^2 T}} \quad \text{or} \quad \eta_t = \frac{1}{1 + \frac{c_2(1-\gamma)t}{\log^2 T}}$$

# Sample complexity of synchronous Q-learning

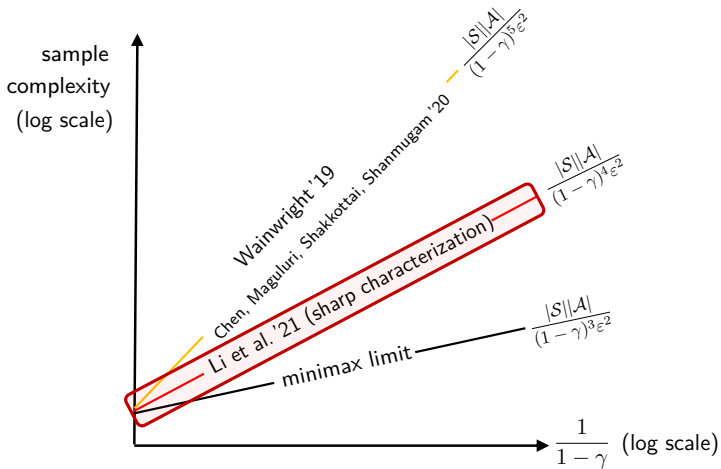
## Theorem 2 (Li, Cai, Chen, Wei, Chi '21)

For any  $0 < \varepsilon \leq 1$ , synchronous Q-learning yields  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  with high prob. and  $\mathbb{E}[\|\hat{Q} - Q^*\|_\infty] \leq \varepsilon$ , with sample size **at most**

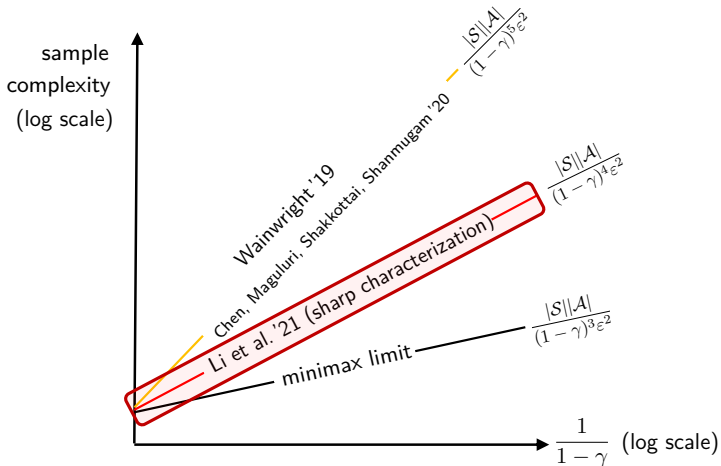
$$\begin{cases} \tilde{O}\left(\frac{|S||\mathcal{A}|}{(1-\gamma)^4 \varepsilon^2}\right) & \text{if } |\mathcal{A}| \geq 2 & (?) \\ \tilde{O}\left(\frac{|S|}{(1-\gamma)^3 \varepsilon^2}\right) & \text{if } |\mathcal{A}| = 1 & (\text{minimax optimal}) \end{cases}$$

other papers	sample complexity
Even-Dar & Mansour '03	$2^{\frac{1}{1-\gamma}} \frac{ S  \mathcal{A} }{(1-\gamma)^4 \varepsilon^2}$
Beck & Srikant '12	$\frac{ S ^2  \mathcal{A} ^2}{(1-\gamma)^5 \varepsilon^2}$
Wainwright '19	$\frac{ S  \mathcal{A} }{(1-\gamma)^5 \varepsilon^2}$
Chen, Maguluri, Shakkottai, Shanmugam '20	$\frac{ S  \mathcal{A} }{(1-\gamma)^5 \varepsilon^2}$

All this requires sample size at least  $\frac{|S||\mathcal{A}|}{(1-\gamma)^4 \epsilon^2}$  ( $|\mathcal{A}| \geq 2$ ) ...



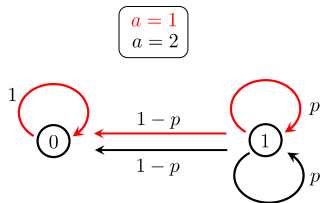
All this requires sample size at least  $\frac{|S||\mathcal{A}|}{(1-\gamma)^4 \epsilon^2}$  ( $|\mathcal{A}| \geq 2$ ) ...



**Question:** Is Q-learning sub-optimal, or is it an analysis artifact?

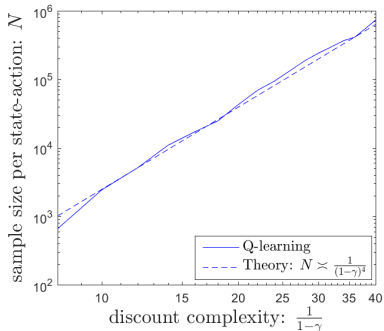
**A numerical example:**  $\frac{|S||A|}{(1-\gamma)^4 \epsilon^2}$  samples seem necessary ...

— *observed in Wainwright '19*



$$p = \frac{4\gamma - 1}{3\gamma}$$

$$r(0, 1) = 0, \quad r(1, 1) = r(1, 2) = 1$$



# Q-learning is NOT minimax optimal

## Theorem 3 (Li, Cai, Chen, Wei, Chi, 2021)

For any  $0 < \varepsilon \leq 1$ , there exists an MDP with  $|\mathcal{A}| \geq 2$  such that to achieve  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ , synchronous Q-learning needs *at least*

$$\tilde{\Omega} \left( \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^4 \varepsilon^2} \right) \text{ samples}$$

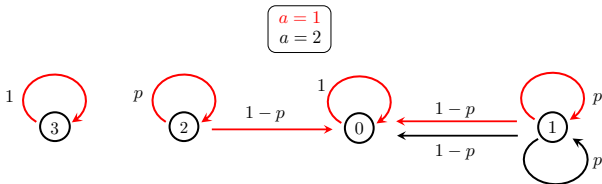
# Q-learning is NOT minimax optimal

## Theorem 3 (Li, Cai, Chen, Wei, Chi, 2021)

For any  $0 < \varepsilon \leq 1$ , there exists an MDP with  $|\mathcal{A}| \geq 2$  such that to achieve  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ , synchronous Q-learning needs *at least*

$$\tilde{\Omega} \left( \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^4 \varepsilon^2} \right) \text{ samples}$$

- Tight **algorithm-dependent** lower bound
- Holds for both constant and rescaled linear learning rates



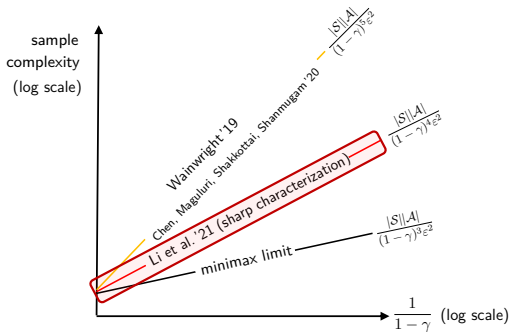


# Q-learning is NOT minimax optimal

## Theorem 3 (Li, Cai, Chen, Wei, Chi, 2021)

For any  $0 < \varepsilon \leq 1$ , there exists an MDP with  $|\mathcal{A}| \geq 2$  such that to achieve  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ , synchronous Q-learning needs *at least*

$$\tilde{\Omega} \left( \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^4 \varepsilon^2} \right) \text{ samples}$$



*Improving sample complexity via **variance reduction***

— *a powerful idea from finite-sum stochastic optimization*

## Variance-reduced Q-learning updates (Wainwright '19)

— inspired by SVRG (Johnson & Zhang '13)

$$Q_t(s, a) = (1 - \eta)Q_{t-1}(s, a) + \eta \left( \mathcal{T}_t(Q_{t-1}) - \underbrace{\mathcal{T}_t(\bar{Q}) + \tilde{\mathcal{T}}(\bar{Q})}_{\text{use } \bar{Q} \text{ to help reduce variability}} \right)(s, a)$$

## Variance-reduced Q-learning updates (Wainwright '19)

— inspired by SVRG (Johnson & Zhang '13)

$$Q_t(s, a) = (1 - \eta)Q_{t-1}(s, a) + \eta \left( \mathcal{T}_t(Q_{t-1}) - \underbrace{\mathcal{T}_t(\bar{Q}) + \tilde{\mathcal{T}}(\bar{Q})}_{\text{use } \bar{Q} \text{ to help reduce variability}} \right)(s, a)$$

- $\bar{Q}$ : some reference Q-estimate
- $\tilde{\mathcal{T}}$ : empirical Bellman operator (using a batch of samples)

$$\mathcal{T}_t(Q)(s, a) = r(s, a) + \gamma \max_{a'} Q(s', a')$$

$$\tilde{\mathcal{T}}(Q)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \tilde{P}(\cdot|s, a)} \left[ \max_{a'} Q(s', a') \right]$$

# An epoch-based stochastic algorithm

---

— inspired by Johnson & Zhang '13

update  $\bar{Q}$  variance-reduced  
Q-learning



**for** each epoch

1. update  $\bar{Q}$  and  $\tilde{\mathcal{T}}(\bar{Q})$  (which stay fixed in the rest of the epoch)
2. run variance-reduced Q-learning updates iteratively

# Sample complexity of variance-reduced Q-learning

## Theorem 4 (Wainwright '19)

For any  $0 < \varepsilon \leq 1$ , sample complexity for **variance-reduced synchronous Q-learning** to yield  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  is at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

- allows for more aggressive learning rates

# Sample complexity of variance-reduced Q-learning

## Theorem 4 (Wainwright '19)

For any  $0 < \varepsilon \leq 1$ , sample complexity for **variance-reduced synchronous Q-learning** to yield  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  is at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

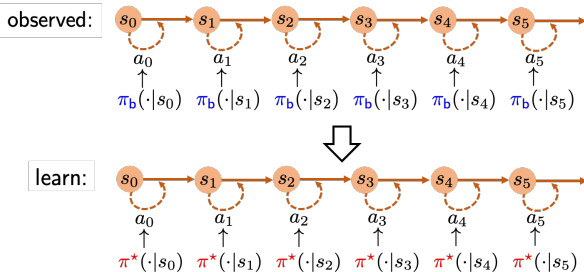
- allows for more aggressive learning rates
- minimax-optimal for  $0 < \varepsilon \leq 1$ 
  - remains suboptimal if  $1 < \varepsilon < \frac{1}{1-\gamma}$

## Model-free RL

1. Basics of Q-learning
2. Synchronous Q-learning and variance reduction (simulator)
3. Asynchronous Q-learning (Markovian data)
4. Q-learning with lower confidence bounds (offline RL)
5. Q-learning with upper confidence bounds (online RL)



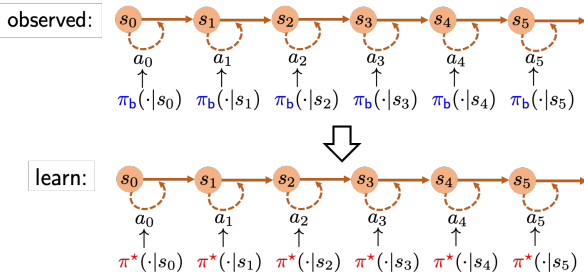
# Markovian samples and behavior policy



**Observed:**  $\underbrace{\{s_t, a_t, r_t\}_{t \geq 0}}$  generated by **behavior policy**  $\pi_b$   
stationary Markovian trajectory

**Goal:** learn optimal value  $V^*$  and  $Q^*$  based on sample trajectory

# Markovian samples and behavior policy



Key quantities of sample trajectory

- minimum state-action occupancy probability (uniform coverage)

$$\mu_{\min} := \min_{\underbrace{(s, a)}_{\text{stationary distribution}}} \mu_{\pi_b}(s, a) \in \left[0, \frac{1}{|\mathcal{S}||\mathcal{A}|}\right]$$

- mixing time:  $t_{\text{mix}}$

# Q-learning on Markovian samples

---



*Chris Watkins*



*Peter Dayan*

$$\underbrace{Q_{t+1}(s_t, a_t) = (1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t)}_{\text{only update } (s_t, a_t)\text{-th entry}}, \quad t \geq 0$$

# Q-learning on Markovian samples

---



*Chris Watkins*

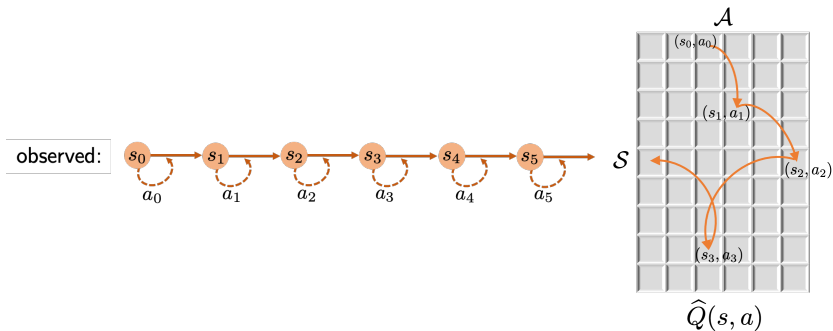


*Peter Dayan*

$$\underbrace{Q_{t+1}(s_t, a_t) = (1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t)}_{\text{only update } (s_t, a_t)\text{-th entry}}, \quad t \geq 0$$

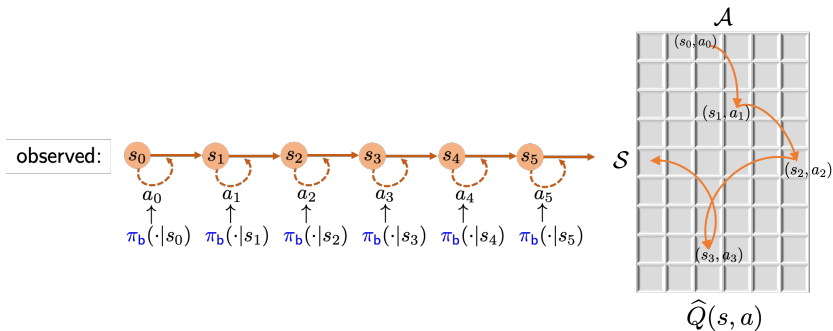
$$\mathcal{T}_t(Q)(s_t, a_t) = r(s_t, a_t) + \gamma \max_{a'} Q(s_{t+1}, a')$$

# Q-learning on Markovian samples



- **asynchronous:** only a single entry is updated each iteration

# Q-learning on Markovian samples



- **asynchronous:** only a single entry is updated each iteration
- **off-policy:** target policy  $\pi^* \neq$  behavior policy  $\pi_b$

# Sample complexity of asynchronous Q-learning

## Theorem 5 (Li, Cai, Chen, Wei, Chi '21)

For any  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , sample complexity of async Q-learning to yield  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  with high prob. (or  $\mathbb{E}[\|\hat{Q} - Q^*\|_\infty] \leq \varepsilon$ ) is at most

$$\frac{1}{\mu_{\min}(1-\gamma)^4\varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)} \quad (\text{up to log factor})$$

# Sample complexity of asynchronous Q-learning

## Theorem 5 (Li, Cai, Chen, Wei, Chi '21)

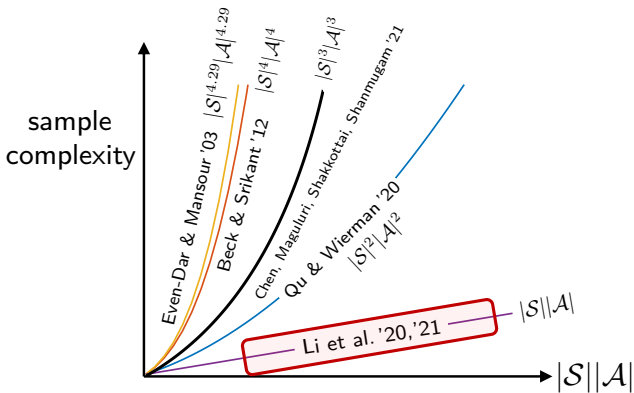
For any  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , sample complexity of async Q-learning to yield  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  with high prob. (or  $\mathbb{E}[\|\hat{Q} - Q^*\|_\infty] \leq \varepsilon$ ) is at most

$$\frac{1}{\mu_{\min}(1-\gamma)^4\varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)} \quad (\text{up to log factor})$$

other papers	sample complexity
Even-Dar, Mansour '03	$\frac{1}{(1-\gamma)^4\varepsilon^2}$
Even-Dar, Mansour '03	$\left(\frac{t_{\text{cover}}^{1+3\omega}}{(1-\gamma)^4\varepsilon^2}\right)^{\frac{1}{\omega}} + \left(\frac{t_{\text{cover}}}{1-\gamma}\right)^{\frac{1}{1-\omega}}, \omega \in (\frac{1}{2}, 1)$
Beck & Srikant '12	$\frac{t_{\text{cover}}^3  S   A }{(1-\gamma)^5 \varepsilon^2}$
Qu & Wierman '20	$\frac{t_{\text{mix}}}{\mu_{\min}^2 (1-\gamma)^5 \varepsilon^2}$
Li, Wei, Chi, Gu, Chen '20	$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$
Chen, Maguluri, Shakkottai, Shanmugam '21	$\frac{1}{\mu_{\min}^3 (1-\gamma)^5 \varepsilon^2} + \text{other-term}(t_{\text{mix}})$



# Linear dependency on $1/\mu_{\min}$



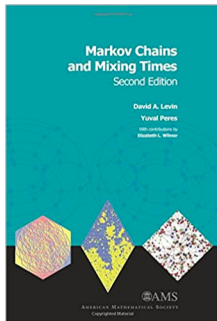
if we take  $\mu_{\min} \asymp \frac{1}{|S||A|}$ ,  $t_{\text{cover}} \asymp \frac{t_{\text{mix}}}{\mu_{\min}}$

# Effect of mixing time on sample complexity

---

$$\frac{1}{\mu_{\min}(1 - \gamma)^4 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1 - \gamma)}$$

- reflects cost taken to reach steady state

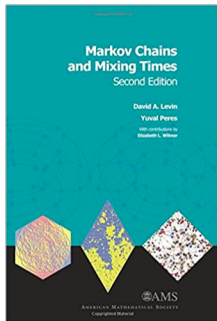


# Effect of mixing time on sample complexity

$$\frac{1}{\mu_{\min}(1-\gamma)^4\varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$

- reflects cost taken to reach steady state
- one-time expense (almost independent of  $\varepsilon$ )
  - it becomes amortized as algorithm runs

— *prior art*:  $\frac{t_{\text{mix}}}{\mu_{\min}^2(1-\gamma)^5\varepsilon^2}$  (Qu & Wierman '20)



## Model-free RL

1. Basics of Q-learning
2. Synchronous Q-learning and variance reduction (simulator)
3. Asynchronous Q-learning (Markovian data)
4. Q-learning with lower confidence bounds (offline RL)
5. Q-learning with upper confidence bounds (online RL)

## Recap: offline RL / batch RL

---

**Historical dataset**  $\mathcal{D} = \{(s^{(i)}, a^{(i)}, s'^{(i)})\}$ :  $N$  independent copies of

$$s \sim \rho^b, \quad a \sim \pi^b(\cdot | s), \quad s' \sim P(\cdot | s, a)$$

for some state distribution  $\rho^b$  and behavior policy  $\pi^b$

## Recap: offline RL / batch RL

**Historical dataset**  $\mathcal{D} = \{(s^{(i)}, a^{(i)}, s'^{(i)})\}$ :  $N$  independent copies of

$$s \sim \rho^b, \quad a \sim \pi^b(\cdot | s), \quad s' \sim P(\cdot | s, a)$$

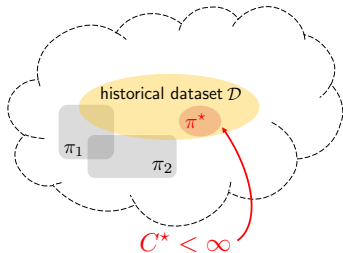
for some state distribution  $\rho^b$  and behavior policy  $\pi^b$

### Single-policy concentrability

$$C^* := \max_{s,a} \frac{d^{\pi^*}(s,a)}{d^{\pi^b}(s,a)} \geq 1$$

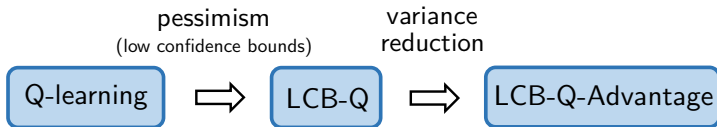
where  $d^\pi$ : occupancy distribution under  $\pi$

- captures **distributional shift**
- allows for partial coverage



*How to design offline model-free algorithms  
with optimal sample efficiency?*

*How to design offline model-free algorithms  
with optimal sample efficiency?*





# LCB-Q: Q-learning with LCB penalty

---

— Shi et al. '22, Yan et al. '22

$$Q_{t+1}(s_t, a_t) \leftarrow \underbrace{(1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t)}_{\text{classical Q-learning}} - \underbrace{\eta_t b_t(s_t, a_t)}_{\text{LCB penalty}}$$

# LCB-Q: Q-learning with LCB penalty

---

— Shi et al. '22, Yan et al. '22

$$Q_{t+1}(s_t, a_t) \leftarrow \underbrace{(1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t)}_{\text{classical Q-learning}} - \underbrace{\eta_t b_t(s_t, a_t)}_{\text{LCB penalty}}$$

- $b_t(s, a)$ : Hoeffding-style confidence bound
- pessimism in the face of uncertainty

# LCB-Q: Q-learning with LCB penalty

— Shi et al. '22, Yan et al. '22

$$Q_{t+1}(s_t, a_t) \leftarrow \underbrace{(1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t)}_{\text{classical Q-learning}} - \underbrace{\eta_t b_t(s_t, a_t)}_{\text{LCB penalty}}$$

- $b_t(s, a)$ : Hoeffding-style confidence bound
- pessimism in the face of uncertainty

sample size:  $\tilde{O}\left(\frac{SC^*}{(1-\gamma)^5 \epsilon^2}\right) \implies$  sub-optimal by a factor of  $\frac{1}{(1-\gamma)^2}$

**Issue:** large variability in stochastic update rules

# Q-learning with LCB and variance reduction

---

— Shi et al. '22, Yan et al. '22

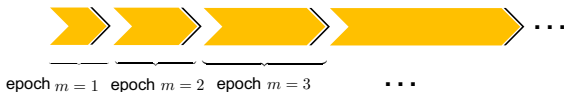
$$Q_{t+1}(s_t, a_t) \leftarrow (1 - \eta_t)Q_t(s_t, a_t) - \eta_t \underbrace{b_t(s_t, a_t)}_{\text{LCB penalty}} + \eta_t \left( \underbrace{\mathcal{T}_t(Q_t) - \mathcal{T}_t(\bar{Q})}_{\text{advantage}} + \underbrace{\hat{\mathcal{T}}(\bar{Q})}_{\text{reference}} \right)(s_t, a_t)$$

# Q-learning with LCB and variance reduction

— Shi et al. '22, Yan et al. '22

$$Q_{t+1}(s_t, a_t) \leftarrow (1 - \eta_t)Q_t(s_t, a_t) - \underbrace{\eta_t b_t(s_t, a_t)}_{\text{LCB penalty}} + \eta_t \left( \underbrace{\mathcal{T}_t(Q_t) - \mathcal{T}_t(\bar{Q})}_{\text{advantage}} + \underbrace{\hat{\mathcal{T}}(\bar{Q})}_{\text{reference}} \right) (s_t, a_t)$$

- incorporates **variance reduction** into LCB-Q

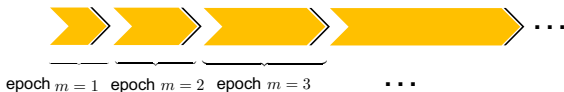


# Q-learning with LCB and variance reduction

— Shi et al. '22, Yan et al. '22

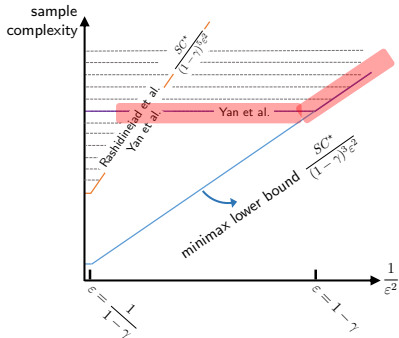
$$Q_{t+1}(s_t, a_t) \leftarrow (1 - \eta_t)Q_t(s_t, a_t) - \underbrace{\eta_t b_t(s_t, a_t)}_{\text{LCB penalty}} + \eta_t \left( \underbrace{\mathcal{T}_t(Q_t) - \mathcal{T}_t(\bar{Q})}_{\text{advantage}} + \underbrace{\hat{\mathcal{T}}(\bar{Q})}_{\text{reference}} \right) (s_t, a_t)$$

- incorporates **variance reduction** into LCB-Q

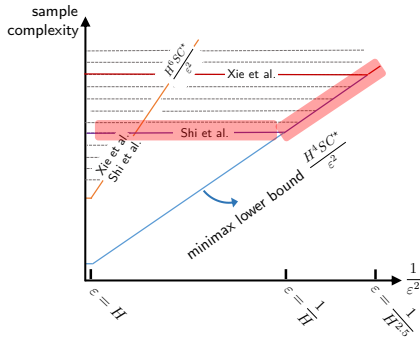


## Theorem 6 (Yan, Li, Chen, Fan '22, Shi, Li, Wei, Chen, Chi '22)

For  $\varepsilon \in (0, 1 - \gamma]$ , LCB-Q-Advantage achieves  $V^*(\rho) - V^{\hat{\pi}}(\rho) \leq \varepsilon$  with optimal sample complexity  $\tilde{O}\left(\frac{SC^*}{(1-\gamma)^3\varepsilon^2}\right)$



infinite-horizon MDPs



finite-horizon MDPs

Model-free offline RL attains sample optimality too!

— with some burn-in cost though ...

## Model-free RL

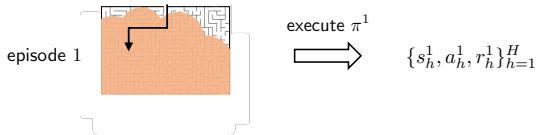
1. Basics of Q-learning
2. Synchronous Q-learning and variance reduction (simulator)
3. Asynchronous Q-learning (Markovian data)
4. Q-learning with lower confidence bounds (offline RL)
5. Q-learning with upper confidence bounds (online RL)



# Online RL: interacting with real environments

---

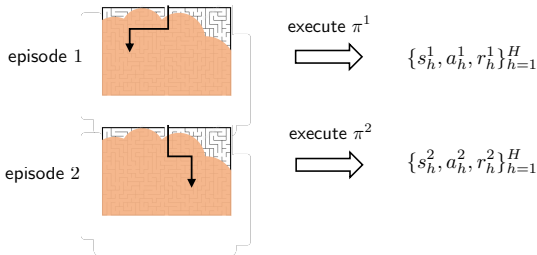
*Sequentially* execute MDP for  $K$  episodes, each consisting of  $H$  steps



# Online RL: interacting with real environments

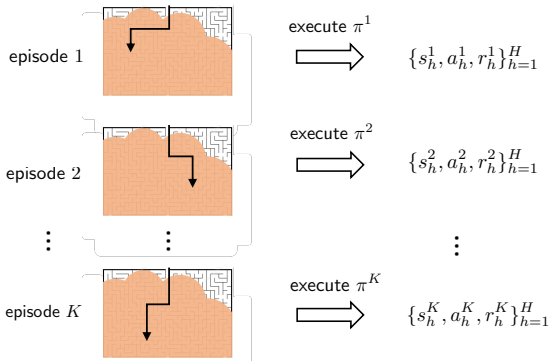
---

*Sequentially* execute MDP for  $K$  episodes, each consisting of  $H$  steps



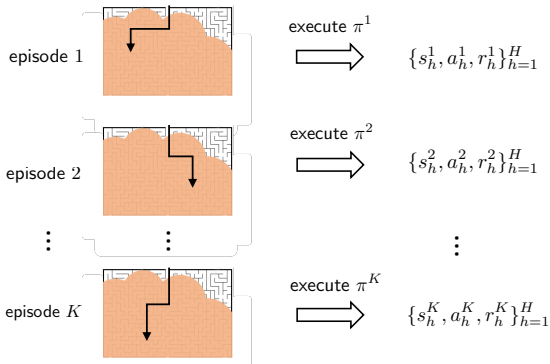
# Online RL: interacting with real environments

*Sequentially* execute MDP for  $K$  episodes, each consisting of  $H$  steps



# Online RL: interacting with real environments

Sequentially execute MDP for  $K$  episodes, each consisting of  $H$  steps  
— *sample size:  $T = KH$*



**exploration** (exploring unknowns) vs. **exploitation** (exploiting learned info)

# Regret: gap between learned policy & optimal policy

---

adversary



learner



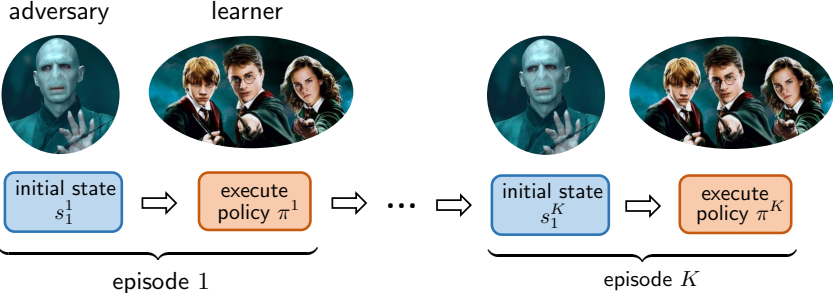
initial state  
 $s_1^1$



execute  
policy  $\pi^1$

episode 1

# Regret: gap between learned policy & optimal policy





## Lower bound

(Domingues et al. '21)

$$\text{Regret}(T) \gtrsim \sqrt{H^2 SAT}$$

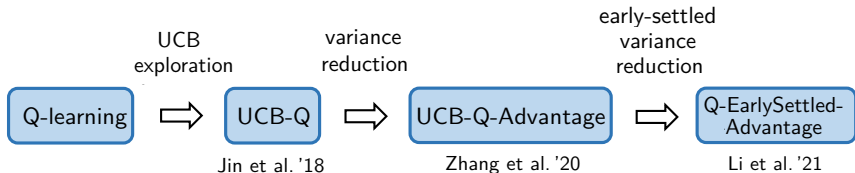
## Existing algorithms

- UCB-VI: Azar et al. '17
- UBEV: Dann et al. '17
- UCB-Q-Hoeffding: Jin et al. '18
- **UCB-Q-Bernstein: Jin et al. '18**
- UCB2-Q-Bernstein: Bai et al. '19
- EULER: Zanette et al. '19
- **UCB-Q-Advantage: Zhang et al. '20**
- UCB-M-Q: Menard et al. '21
- **Q-EarlySettled-Advantage: Li et al. '21**



*Which model-free algorithms are sample-efficient for online RL?*

*Which model-free algorithms are sample-efficient for online RL?*



# Q-learning with UCB exploration (Jin et al., 2018)

---

$$Q_h(s_h, a_h) \leftarrow \underbrace{(1 - \eta_k)Q_h(s_h, a_h) + \eta_k \mathcal{T}_k(Q_{h+1})(s_h, a_h)}_{\text{classical Q-learning}} + \eta_k \underbrace{b_h(s_h, a_h)}_{\text{exploration bonus}}$$

# Q-learning with UCB exploration (Jin et al., 2018)

---

$$Q_h(s_h, a_h) \leftarrow \underbrace{(1 - \eta_k)Q_h(s_h, a_h) + \eta_k \mathcal{T}_k(Q_{h+1})(s_h, a_h)}_{\text{classical Q-learning}} + \eta_k \underbrace{b_h(s_h, a_h)}_{\text{exploration bonus}}$$

- $b_h(s, a)$ : upper confidence bound; encourage exploration  
— *optimism in the face of uncertainty*
- inspired by UCB bandit algorithm (Lai, Robbins '85)

# Q-learning with UCB exploration (Jin et al., 2018)

$$Q_h(s_h, a_h) \leftarrow \underbrace{(1 - \eta_k)Q_h(s_h, a_h) + \eta_k \mathcal{T}_k(Q_{h+1})(s_h, a_h)}_{\text{classical Q-learning}} + \eta_k \underbrace{b_h(s_h, a_h)}_{\text{exploration bonus}}$$

- $b_h(s, a)$ : upper confidence bound; encourage exploration  
— *optimism in the face of uncertainty*
- inspired by UCB bandit algorithm (Lai, Robbins '85)

Regret( $T$ )  $\lesssim \sqrt{H^3 S A T}$   $\implies$  sub-optimal by a factor of  $\sqrt{H}$

## Q-learning with UCB exploration (Jin et al., 2018)

$$Q_h(s_h, a_h) \leftarrow \underbrace{(1 - \eta_k)Q_h(s_h, a_h) + \eta_k \mathcal{T}_k(Q_{h+1})(s_h, a_h)}_{\text{classical Q-learning}} + \eta_k \underbrace{b_h(s_h, a_h)}_{\text{exploration bonus}}$$

- $b_h(s, a)$ : upper confidence bound; encourage exploration  
— *optimism in the face of uncertainty*
- inspired by UCB bandit algorithm (Lai, Robbins '85)

Regret( $T$ )  $\lesssim \sqrt{H^3 S A T}$   $\implies$  sub-optimal by a factor of  $\sqrt{H}$

**Issue:** large variability in stochastic update rules

# UCB Q-learning with UCB and variance reduction

---

Incorporates **variance reduction** into UCB-Q: — *Zhang, Zhou, Ji '20*

- asymptotically regret-optimal

# UCB Q-learning with UCB and variance reduction

---

Incorporates **variance reduction** into UCB-Q: — Zhang, Zhou, Ji '20

- asymptotically regret-optimal
- **Issue:** high burn-in cost  $O(S^6 A^4 H^{28})$



# UCB Q-learning with UCB and variance reduction

---

Incorporates **variance reduction** into UCB-Q: — *Zhang, Zhou, Ji '20*

- asymptotically regret-optimal
- **Issue:** high burn-in cost  $O(S^6 A^4 H^{28})$

One additional idea: early settlement of reference updates — *Li, Shi, Chen, Chi '23*

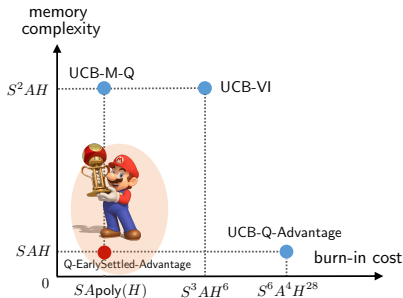
# UCB Q-learning with UCB and variance reduction

Incorporates **variance reduction** into UCB-Q: — Zhang, Zhou, Ji '20

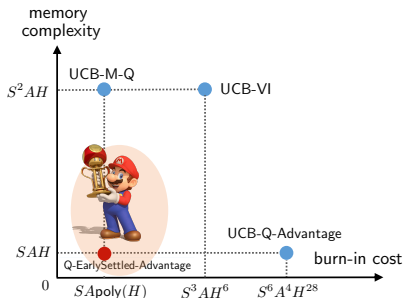
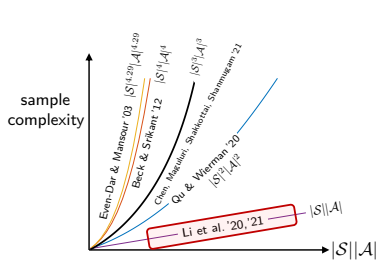
- asymptotically regret-optimal
- **Issue:** high burn-in cost  $O(S^6 A^4 H^{28})$

One additional idea: early settlement of reference updates — Li, Shi, Chen, Chi '23

- regret-optimal w/ near-minimal burn-in cost in  $S$  and  $A$
- memory-efficient  $O(SAH)$
- computationally efficient: runtime  $O(T)$



# Summary of this part



Model-free RL can achieve memory efficiency, computational efficiency, and sample efficiency at once!  
— *with some burn-in cost though*

# Reference I

---

- "*Model-based multi-agent RL in zero-sum Markov games with near-optimal sample complexity*," K. Zhang, S. Kakade, T. Basar, L. Yang, *NeurIPS*, 2020
- "*When can we learn general-sum Markov games with a large number of players sample-efficiently?*" Z. Song, S. Mei, Y. Bai, *ICLR* 2022
- "*V-learning: A simple, efficient, decentralized algorithm for multiagent RL*," C. Jin, Q. Liu, Y. Wang, T. Yu, 2021
- "*Minimax-optimal multi-agent RL in markov games with a generative model*," G. Li, Y. Chi, Y. Wei, Y. Chen, *NeurIPS*, 2022
- "*The complexity of Markov equilibrium in stochastic games*," C. Daskalakis, N. Golowich, K. Zhang, *COLT*, 2023
- "*A stochastic approximation method*," H. Robbins, S. Monro, *Annals of mathematical statistics*, 1951

## Reference II

---

- "*Robust stochastic approximation approach to stochastic programming*," A. Nemirovski, A. Juditsky, G. Lan, A. Shapiro, *SIAM Journal on optimization*, 2009
- "*Learning from delayed rewards*," C. Watkins, 1989
- "*Q-learning*," C. Watkins, P. Dayan, *Machine learning*, 1992
- "*Learning to predict by the methods of temporal differences*," R. Sutton, *Machine learning*, 1988
- "*Analysis of temporal-difference learning with function approximation*," B. van Roy, J. Tsitsiklis, *IEEE transactions on automatic control*, 1997
- "*Learning Rates for Q-learning*," E. Even-Dar, Y. Mansour, *Journal of machine learning Research*, 2003
- "*The asymptotic convergence-rate of Q-learning*," C. Szepesvari, *NeurIPS*, 1998

## Reference III

---

- "Stochastic approximation with cone-contractive operators: Sharp  $\ell_\infty$  bounds for Q-learning," M. Wainwright, arXiv:1905.06265, 2019
- "Is Q-Learning minimax optimal? A tight sample complexity analysis," G. Li, Y. Wei, Y. Chi, Y. Chen, accepted to *Operations Research*, 2023
- "Accelerating stochastic gradient descent using predictive variance reduction," R. Johnson, T. Zhang, *NeurIPS*, 2013
- "Variance-reduced Q-learning is minimax optimal," M. Wainwright, arXiv:1906.04697, 2019
- "Asynchronous stochastic approximation and Q-learning," J. Tsitsiklis, *Machine learning*, 1994
- "On the convergence of stochastic iterative dynamic programming algorithms," T. Jaakkola, M. Jordan, S. Singh, *Neural computation*, 1994

## Reference IV

---

- "*Error bounds for constant step-size Q-learning*," C. Beck, R. Srikant, *Systems and control letters*, 2012
- "*Sample complexity of asynchronous Q-learning: sharper analysis and variance reduction*," G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, *NeurIPS* 2020
- "*Finite-time analysis of asynchronous stochastic approximation and Q-learning*," G. Qu, A. Wierman, *COLT* 2020.
- "*Pessimistic Q-learning for offline reinforcement learning: Towards optimal sample complexity*," L. Shi, G. Li, Y. Wei, Y. Chen, Y. Chi, *ICML* 2022.
- "*The efficacy of pessimism in asynchronous Q-learning*," Y. Yan, G. Li, Y. Chen, J. Fan, arXiv:2203.07368, 2022.
- "*Asymptotically efficient adaptive allocation rules*," T. L. Lai, H. Robbins, *Advances in applied mathematics*, vol. 6, no. 1, 1985.

## Reference V

---

- "*Is Q-learning provably efficient?*" C. Jin, Z. Allen-Zhu, S. Bubeck, and M. Jordan, *NeurIPS* 2018.
- "*Almost optimal model-free reinforcement learning via reference-advantage decomposition,*" Z. Zhang, Y. Zhou, X. Ji, *NeurIPS* 2020.
- "*Breaking the sample complexity barrier to regret-optimal model-free reinforcement learning,*" G. Li, L. Shi, Y. Chen, Y. Chi, *Information and Inference: A Journal of the IMA*, 2023.